

FYS-STK4155 - Project 1

Geir Tore Ulvik, Idun Klvstad

October 3, 2019

Abstract

This project looks at three different regression methods, Ordinary Least Squares, Lasso and Ridge, used on a known function together with the bootstrap resampling method. It compares the three methods by looking at the Mean Squared error and the R^2 score. The results show that the OLS method performs the best when trying to estimate Franke's function.

1 Introduction

This report focuses on regression analysis and resampling methods. Three different regression methods will be combined with the bootstrap method and compared. The regression models that will be looked at are the Ordinary Least Squares (OLS) method, Ridge regression and Lasso regression.

The first thing we will study is how to fit polynomials of different order to a specific two-dimensional function called Franke's function. Thereafter the same methods will be used, trying to reproduce digital terrain data.

The report is structured so that there is first given a bit of theory about the different methods, then we will present how we chose to implement those methods and why, before we show and discuss the results of using the different methods and finally makes a conclusion of what we think are each methods strong and weak sides.

2 Theory

2.1 Franke's function

Franke's function is a weighted sum of four exponentials and has been widely used when testing various interpolation and fitting algorithms. The function is given by

$$f(x, y) = \frac{3}{4} \exp \left(-\frac{(9x-2)^2}{4} - \frac{(9y-2)^2}{4} \right) + \frac{3}{4} \exp \left(-\frac{(9x+1)^2}{49} - \frac{(9y+1)^2}{10} \right) \\ + \frac{1}{2} \exp \left(-\frac{(9x-7)^2}{4} - \frac{(9y-3)^2}{4} \right) - \frac{1}{5} \exp \left(-(9x-4)^2 - (9y-7)^2 \right).$$

. [1]

2.2 Linear methods for regression

A linear model assumes that the regression function is linear in the inputs X_1, \dots, X_p , and can give a polynomial representation using basic expansions such as $X_2 = X_1^2, X_3 = X_1^3$.

If we have an input vector as described above and want to predict a real-valued output Y , the linear model reads as

$$f(X) = \beta_0 + \sum_{j=1}^p X_j \beta_j \quad (1)$$

[2]

For all three regression methods that we will look at in this report, there is a set of training data from which we want to estimate the parameter β .

2.2.1 Ordinary Least Squares method

The Ordinary Least Squares (OLS) method pick the coefficients $\beta = (\beta_0, \beta_1, \dots, \beta_p)^T$ to minimize the residual sum of squares given by

$$RSS(\beta) = \sum_{i=1}^N (y_i - f(x_i))^2 \quad (2)$$

$$= \sum_{i=1}^N (y_i - \beta_0 - \sum_{j=1}^p X_{ij} \beta_j)^2 \quad (3)$$

$$(4)$$

This can be rewritten with matrix notation

$$RSS(\beta) = (\hat{y} - \vec{X}\beta)^T (\hat{y} - \vec{X}\beta) \quad (5)$$

[2] As this is explained in detail in the book as well as in the lecture notes [3], we will not derive it here.

As we want to minimize the residual sum of squares, we require

$$\frac{\partial RSS(\hat{\beta})}{\partial \hat{\beta}} = 0 = \hat{X}^T(\hat{y} - \hat{X}\hat{\beta}) \quad (6)$$

This can be rewritten as

$$\hat{X}^T \hat{y} = \hat{X}^T \hat{X} \hat{\beta} \quad (7)$$

We now have to assume that the matrix $\hat{X}^T \hat{X}$ is invertible to get the solution

$$\hat{\beta} = (\hat{X}^T \hat{X})^{-1} \hat{X}^T \hat{y} \quad (8)$$

[3] [2]

2.2.2 Ridge regression

The OLS will not have a solution if the design matrix \hat{X} is singular or near singular, as it depends on $\hat{X}^T \hat{X}$ being invertible. [3]

Ridge regression is known as a shrinkage method. It shrinks the regression coefficients $\hat{\beta}$ by imposing a penalty on their size [2]. We are in other words adding a diagonal component to the matrix to invert. From OLS we therefore make the following change

$$\hat{X}^T \hat{X} \rightarrow \hat{X}^T \hat{X} + \lambda \hat{I} \quad (9)$$

where \hat{I} is the identity matrix. [3]. Thus we get

$$RSS(\lambda) = (\hat{y} - \hat{X}\hat{\beta})^T (\hat{y} - \hat{X}\hat{\beta}) - \lambda \hat{\beta}^T \hat{\beta} \quad (10)$$

and

$$\hat{\beta}^{ridge} = (\hat{X}^T \hat{X} + \lambda \hat{I})^{-1} \hat{X}^T \hat{y} \quad (11)$$

$\lambda > 0$ is a complexity parameter that controls the amount of shrinkage. The larger value of λ the greater amount of shrinkage. [2]

This can also be written in a non-vector format which makes explicit the size constraint on the parameters. It also makes it easier to see the difference from Lasso regression:

$$\hat{\beta}^{ridge} = \underset{\beta}{\operatorname{argmin}} \sum_{i=1}^N \left(y_i - \beta_0 - \sum_{j=1}^p x_{ij} \beta_j \right)^2 \quad (12)$$

subject to

$$\sum_{j=1}^p \beta_j^2 \leq t \quad (13)$$

There is a one to one correspondence between λ and t . [2]

2.2.3 Lasso regression

Like ridge, lasso is a shrinkage method. The main difference from ridge is that lasso has no closed form expression. [2]

Similarly to ridge, the lasso estimate is given by

$$\hat{\beta}^{lasso} = \underset{\beta}{\operatorname{argmin}} \sum_{i=1}^N \left(y_i - \beta_0 - \sum_{j=1}^p x_{ij} \beta_j \right)^2 \quad (14)$$

subject to

$$\sum_{j=1}^p |\beta_j| \leq t \quad (15)$$

As you can see, the ridge penalty $\sum_1^p \beta_j^2$ is replaced by a lasso penalty given by $\sum_1^p |\beta_j|$. This makes the solution nonlinear in the y_i and is the reason that the lasso solution don't have a closed form expression.

Lasso does a kind of continous subset selection because of the nature of the constraint. The reason for that is that making t sufficiently small, will cause some of the coefficients to go to zero. [2]

2.3 Error estimates

To evaluate the different regression methods we have two main measures that can be used. In both equation below \tilde{y}_i is the predicted value of the i-th sample and that y_i is the corresponding true value. [1] The first is the Mean Squared Error (MSE) which is given by

$$MSE(\hat{y}, \tilde{y}) = \frac{1}{n} \sum_{i=0}^{n-1} (y_i - \tilde{y}_i)^2 \quad (16)$$

The other is the R^2 score which is defined by

$$R^2(\hat{y}, \tilde{y}) = 1 - \frac{\sum_{i=0}^{n-1} (y_i - \tilde{y}_i)^2}{\sum_{i=0}^{n-1} (y_i - \bar{y})^2}, \quad (17)$$

where we have defined the mean value of \hat{y} as

$$\bar{y} = \frac{1}{n} \sum_{i=0}^{n-1} y_i. \quad (18)$$

[1]

2.4 Resampling methods

A resampling method is a tool that involves repeatedly drawing samples from a training data set and refitting a model of interest on each sample, in order to obtain additional information about the fitted model. The reason for using a resampling method, is to obtain information that would not be available from fitting the model only once using the original training sample. [3]

2.4.1 Bootstrap

Bootstrap is a resampling method, and is widely used. The bootstrap method is based on the fact that $\hat{\Theta} = \hat{\Theta}$ is a random variable, because it is a function of random variables. Therefore it has a pdf, $p(\vec{t})$. The aim of the bootstrap is to estimate $p(\vec{t})$ by the relative frequency of $\hat{\Theta}$ [3]

The bootstrap method works with four different steps:

1. Draw with replacement n numbers for the observed variables $\hat{x} = (x_1, x_2, \dots, x_n)$.
2. Define a vector \hat{x}^* containing the values which were drawn from \hat{x} .
3. Using the vector \hat{x}^* compute $\hat{\Theta}^*$ by evaluating $\hat{\Theta}$ under the observations \hat{x}^* .
4. Repeat the process k times.

[3]

3 Implementation of methods

All code can be found on github, using the following link:
<https://github.com/geirtul/fys-stk4155/tree/master/project1/src>

There is one class for each regression method, with corresponding filenames. Mean Squared Error and R2 Score are implemented as class methods.

In all the methods, Franke's function is defined for $x, y \in [0, 1]$.

3.1 Ordinary Least Squares

The class ols.py is made to make a prediction of Franke's function based on the OLS method. It makes x, y as arrays with uniformly distributed random numbers between 0 and 1.

It uses scikit learn to make a polynomial of a given degree.

The following lines calculates $\hat{\beta}$ as described in section 2.2.1

```
def mean_squared_error(self, x, y):  
    """  
    Evaluate the mean squared error of the output generated  
    by Ordinary Least Squares regressions.
```

Then the following lines makes a predicted outcome:

```
y_predict = self.make_prediction(x)  
mse = np.mean(np.square(y - y_predict))  
  
return mse
```

3.2 Ridge

As stated in section 2.2.2, there is not a huge difference from OLS to ridge regression.

The code for calculating $\hat{\beta}$ from the math shown in the section looks like this:

```
return y_predict  
  
def mean_squared_error(self, x, y):  
    """  
    Evaluate the mean squared error of the output generated  
    by Ordinary Least Squares regressions.
```

As we only change β , the code for finding the estimated function is the same, only with the new beta.

3.3 Lasso

We chose to use scikit learn for the lasso regression. The code should do what is described in 2.2.3.

The coefficients is therefore calculated like this

```
return y_predict

def mean_squared_error(self, x, y):
    """
    Evaluate the mean squared error of the output generated
    by regressions.
```

and the predicted function like this

```
y_predict = self.lasso_object.predict(x)
mse = np.mean(np.square(y - y_predict))
```

3.4 Bootstrapping

We chose to use bootstrap because of a number of advantages listed in the lecture notes. [3]

The data is first split into training and test sets. Then the bootstrap method described in ?? is implemented as follows

```
# Analysis
# Need y_test[:, None] for proper broadcasting of arrays.
error_out = np.mean(np.mean((y_test[:, None] - y_pred) **
                             2,
                             axis=1, keepdims=True))
var_error_out = np.var(np.mean((y_test[:, None] - y_pred)
                                ** 2,
                                axis=1, keepdims=True))
```

[4]

4 Results

4.1 Ordinary Least Squares Regression on Franke Function data

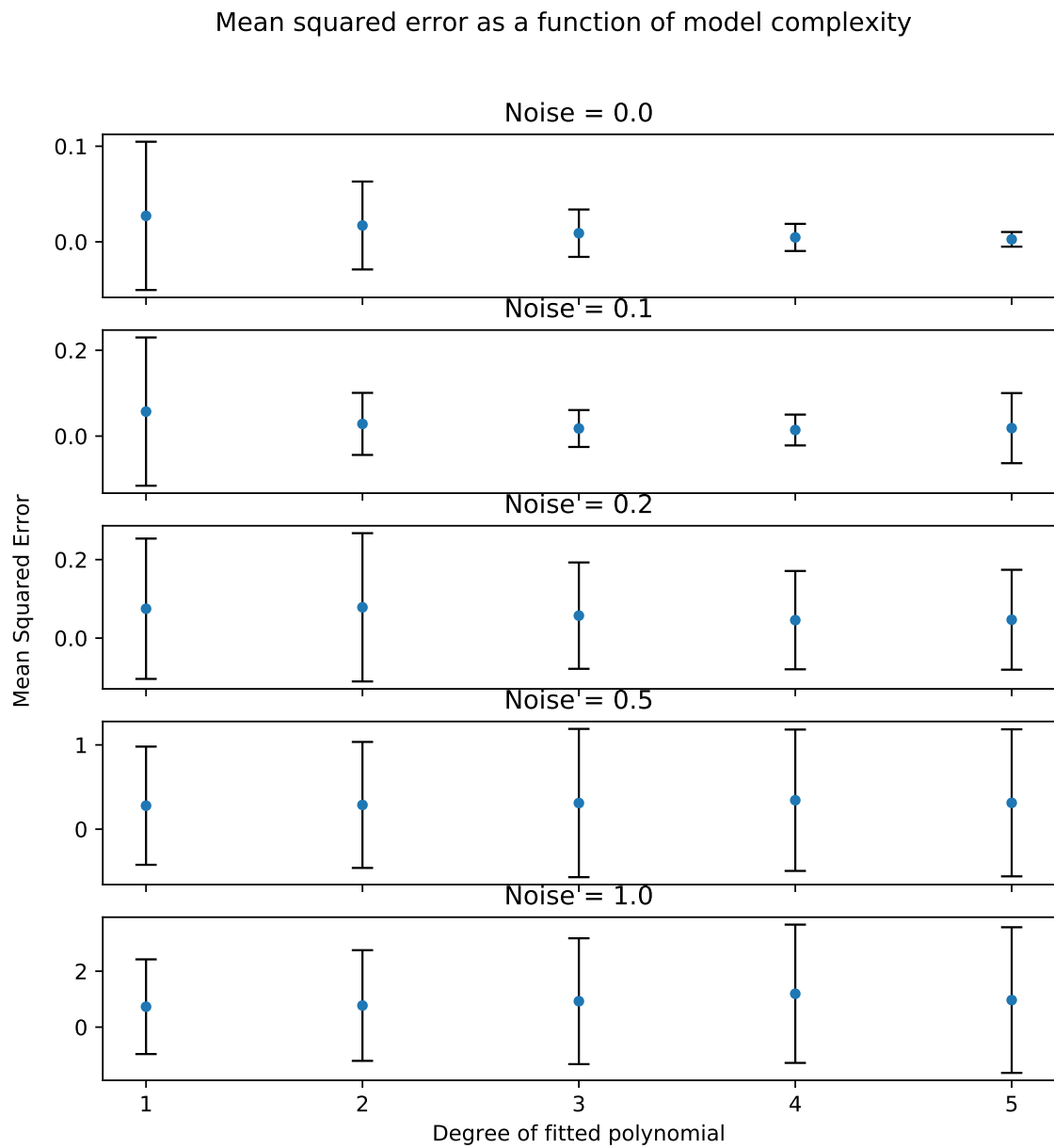


Figure 1: Plot of Mean Squared Error as a function of model complexity (degree of the fitted polynomial) when fitting with Ordinary Least Squares regression. Each plot from top to bottom shows results for a different degree of stochastic noise added to the training data.

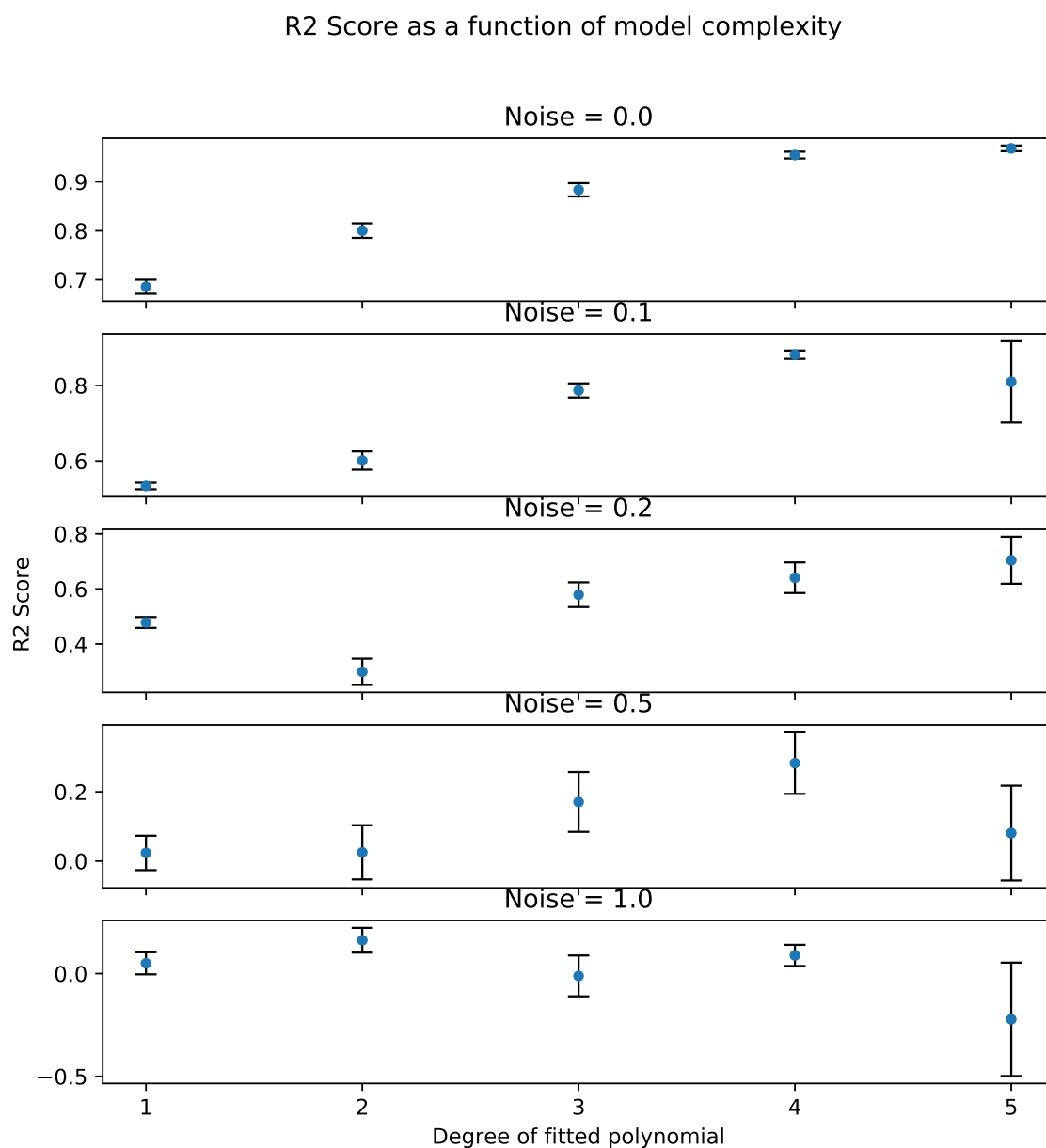


Figure 2: Plot of R2 score as a function of model complexity (degree of the fitted polynomial) when fitting with Ordinary Least Squares regression. Each plot from top to bottom shows results for a different degree of stochastic noise added to the training data.

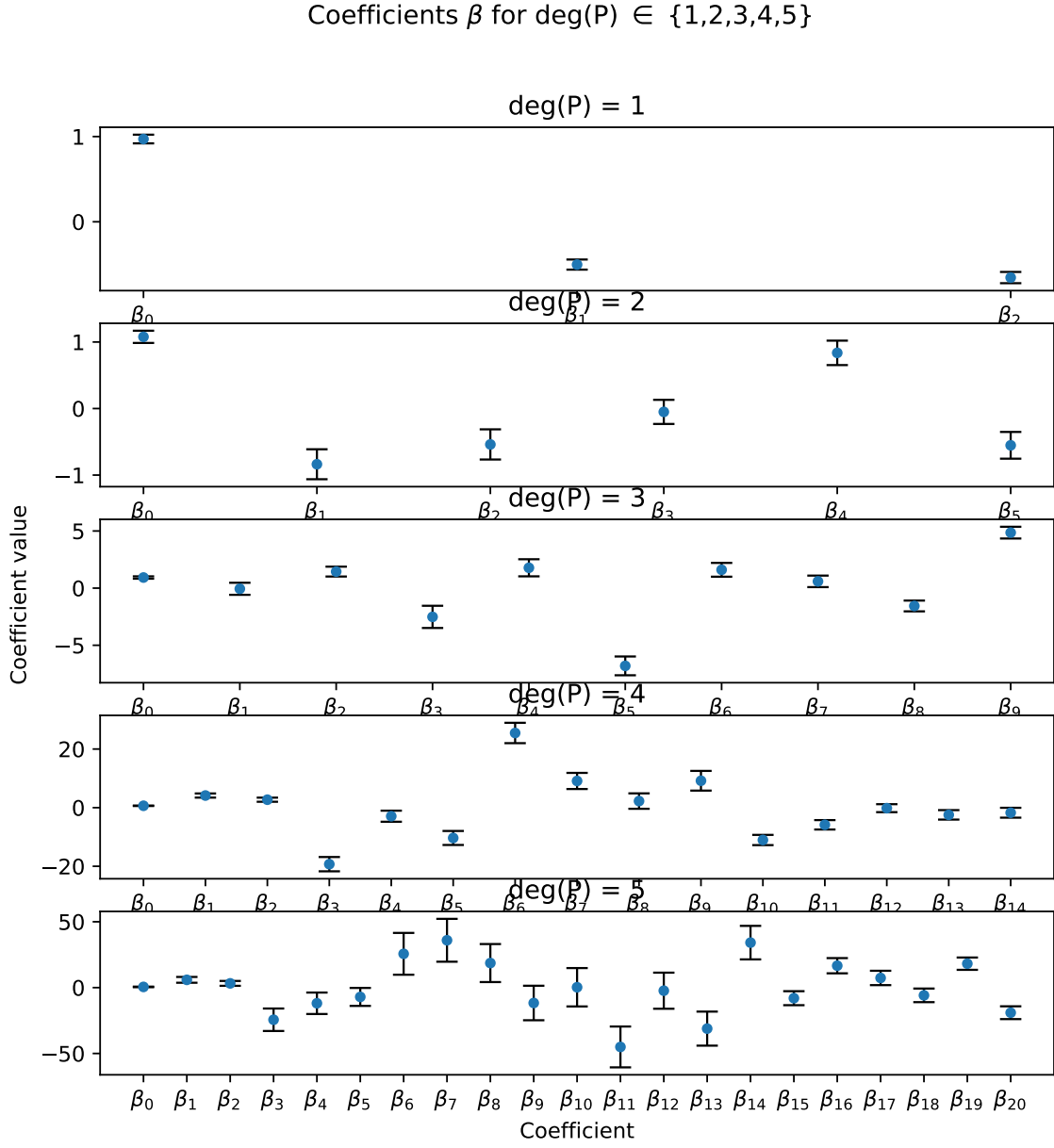


Figure 3: Coefficients for each individual term of the fitted polynomial, after regression. Each coefficient is shown with it's confidence interval (95%), calculated from the variance obtained through bootstrapping.

Coefficients β for $\deg(P) = 5$, with noise

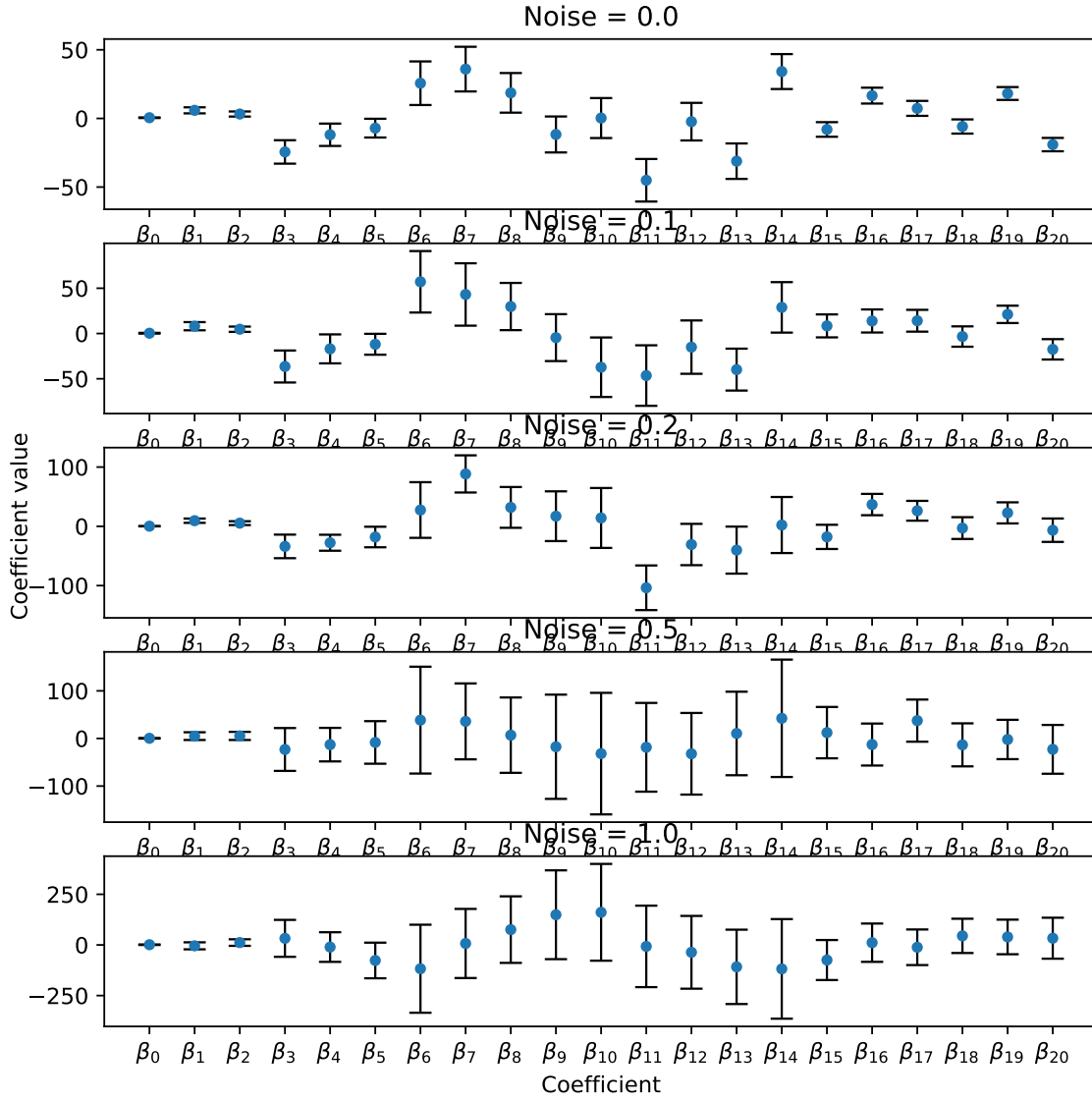


Figure 4: Coefficients for the best polynomial fit based on MSE values. Each plot from top to bottom shows the coefficients with their confidence interval (95%), calculated from the variance obtained through bootstrapping, for varying degrees of stochastic noise added to the training data.

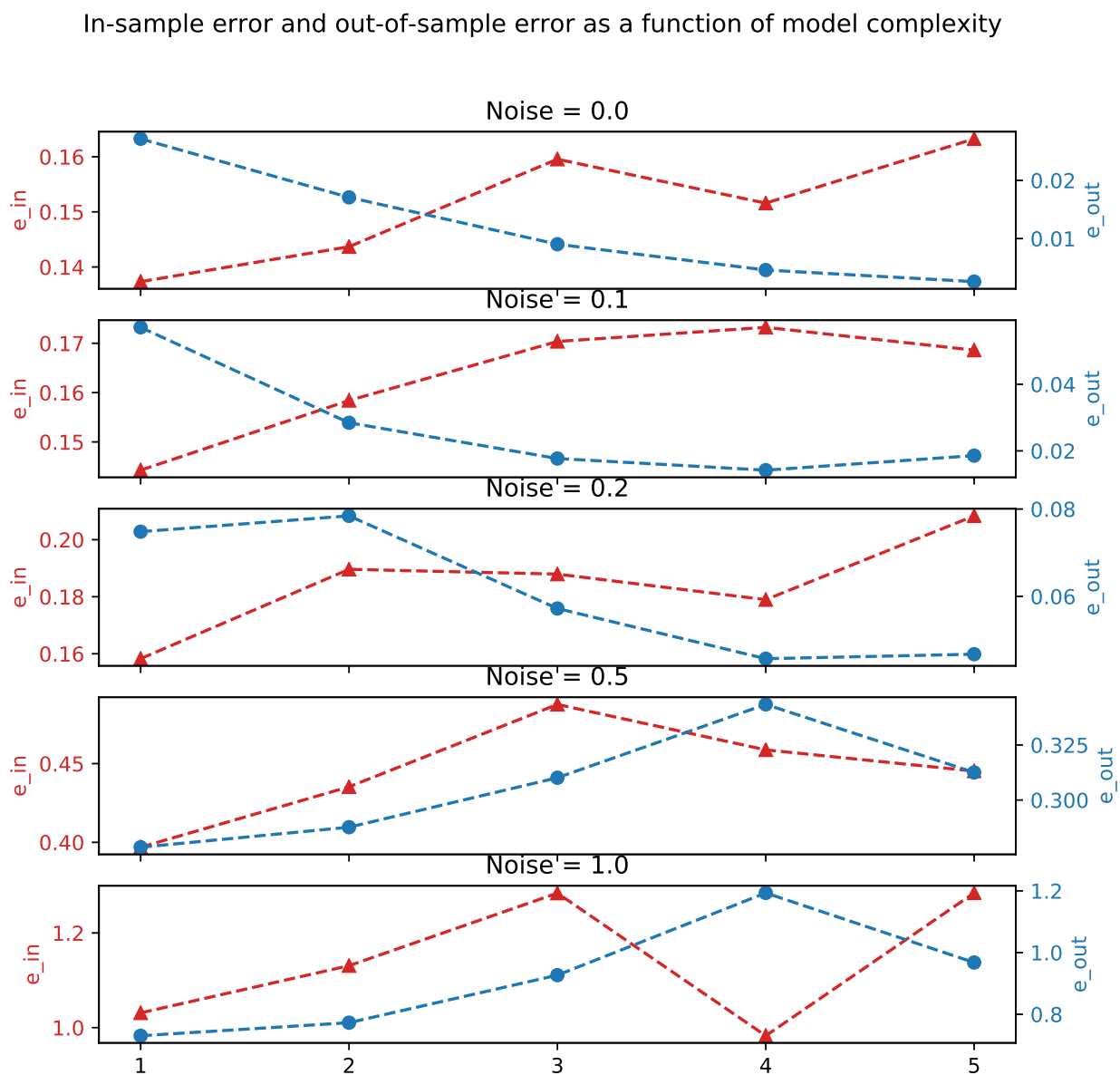


Figure 5: In-sample and out-of-sample errors plotted as a function of model complexity (degree of fitted polynomial).

5 Discussion

The plot in figure ?? show the MSE declining with model complexity, as would be expected for the Franke Function dataset, since it actually includes a polynomial of degree 5. As can be seen from the visualization in ?? this is indeed the case, as the surface imitates that of the dataset quite well, but not perfectly. For the real-world terrain data, the MSE follows a similar trend, shown in figure ??. The data obtained in this report are insufficient to conclude which model is the best suited to use for the terrain data. Further investigation using bootstrap or k-fold cross validation, and analysing the variance and bias of the different models would likely improve the possibility of a conclusion considerably.

6 Conclusion

References

- [1] University of Oslo Department of Physics. Project 1 on machine learning. <https://compphysics.github.io/MachineLearning/doc/Projects/2018/Project1/html/Project1-bs.html>, 2018.
- [2] Gareth James, Daniela Witten, Trevor Hastie, and Robert Tibshirani. *An Introduction to Statistical Learning: With Applications in R*. Springer Publishing Company, Incorporated, 2014.
- [3] M Hjorth-Jensen. Data analysis and machine learning: Linear regression and more advanced regression analysis - lecture notes. https://compphysics.github.io/MachineLearning/doc/pub/Regression/html/._Regression-bs000.html, 2018.
- [4] Bendik Samseth. Piazza post on bias and variance. <https://piazza.com/class/ji78s1cdul39a?cid=59>, 2018.